

Non-coding RNAs: how to find and make sense of them

Judith Zaugg

Luscombe group at EBI

zaugg@ebi.ac.uk

Outline of this lecture

- Introduction about non-coding RNAs
- Methods to identify non-coding RNAs
- Methods to assign biological/molecular functions to non-coding RNAs

INTRODUCTION

“Although less than 2% of a mammalian genome codes for protein, studies consistently show that half or even more of the genome is transcribed”

Long noncoding RNAs: the search for function

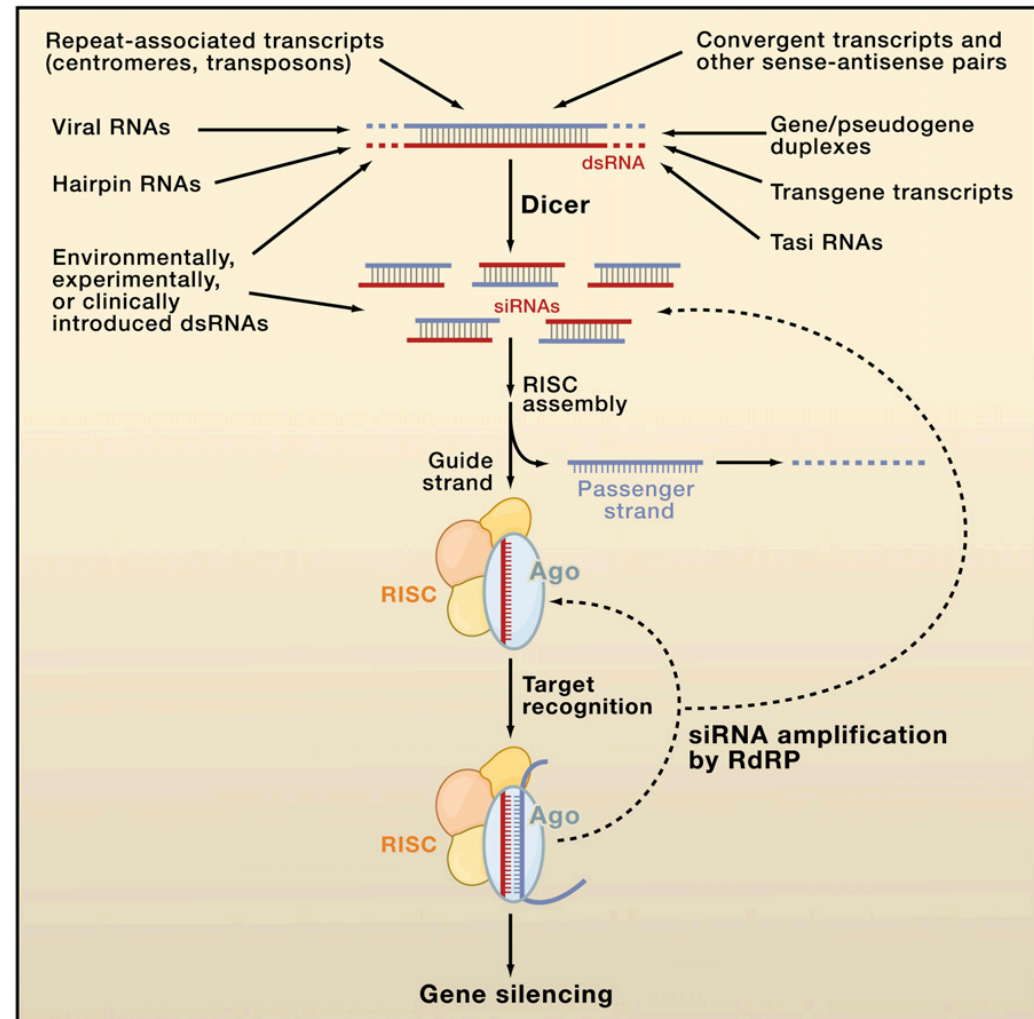
Monya Baker

What are non-coding RNAs

- RNAs that do not encode proteins
- “Textbook knowledge” non-coding (nc)RNA
 - *transfer* (t)RNAs: transfer amino acids
 - *ribosomal* (r)RNAs: constitute the catalytic center of the ribosome
 - *small nucleolar* (sn)RNAs: involved in nuclear processes such as splicing or rRNA processing
- small regulatory RNAs like *small interfering* (si)RNAs
 - gene silencing through double-stranded RNA
- long non-coding RNAs: focus of this presentation

small regulatory RNAs

- just a short overview, not discussed here in detail



Leading Edge
Review

Origins and Mechanisms of miRNAs and siRNAs

Richard W. Carthew^{1,*} and Erik J. Sontheimer^{1,*}
¹Department of Biochemistry, Molecular Biology, and Cell Biology, Northwestern University, 2205 Tech Drive, Evanston, IL 60208-3500, USA
*Correspondence: r-carthew@northwestern.edu (R.W.C.), erik@northwestern.edu (E.J.S.)
DOI 10.1016/j.cell.2009.01.035

long non-coding RNAs (>200bp)

Recent reviews of long non-coding RNAs have focused on:

- function of ncRNAs and mainly describe examples
- how ncRNAs are transcribed

Leading Edge
Review

APPLICATIONS OF NEXT-GENERATION SEQUENCING

The complex eukaryotic transcriptome:
unexpected pervasive transcription
and novel small RNAs

Alain Jacquier

Cell

Evolution and Functions of Long Noncoding RNAs

Chris P. Ponting,^{1,*} Peter L. Oliver,¹ and Wolf Reik²

Long non-coding RNAs: insights into functions

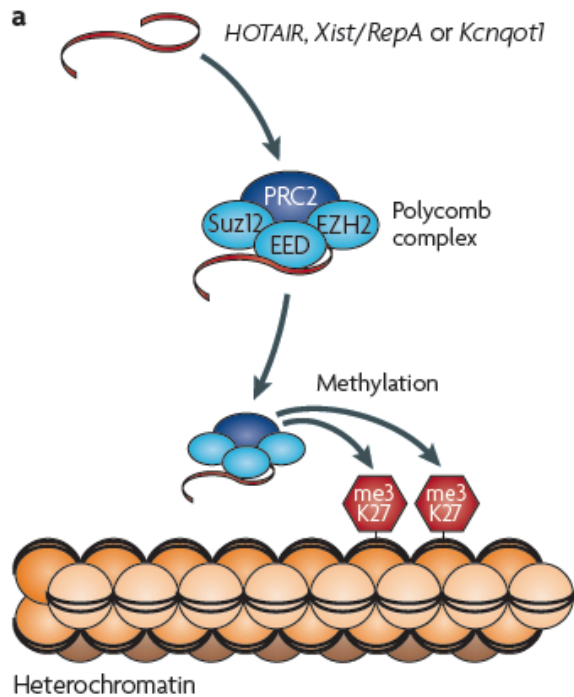
Tim R. Mercer, Marcel E. Dinger and John S. Mattick

Long noncoding RNAs: the search for function

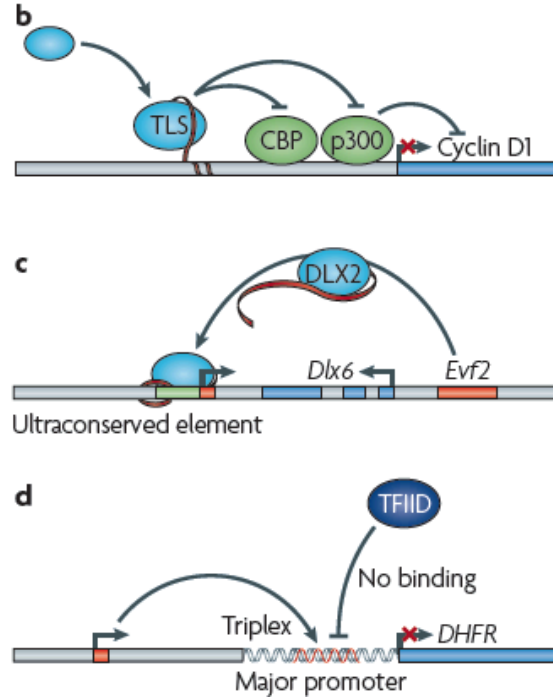
Monya Baker

Examples of how ncRNAs can function

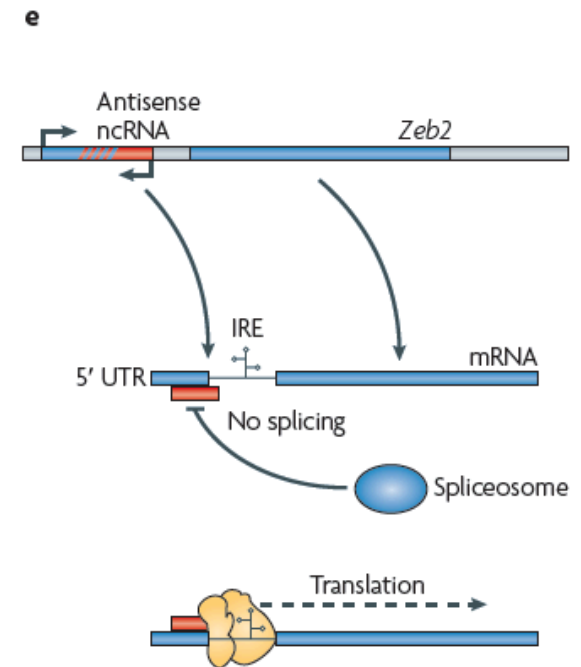
Chromatin remodelling



Transcriptional control

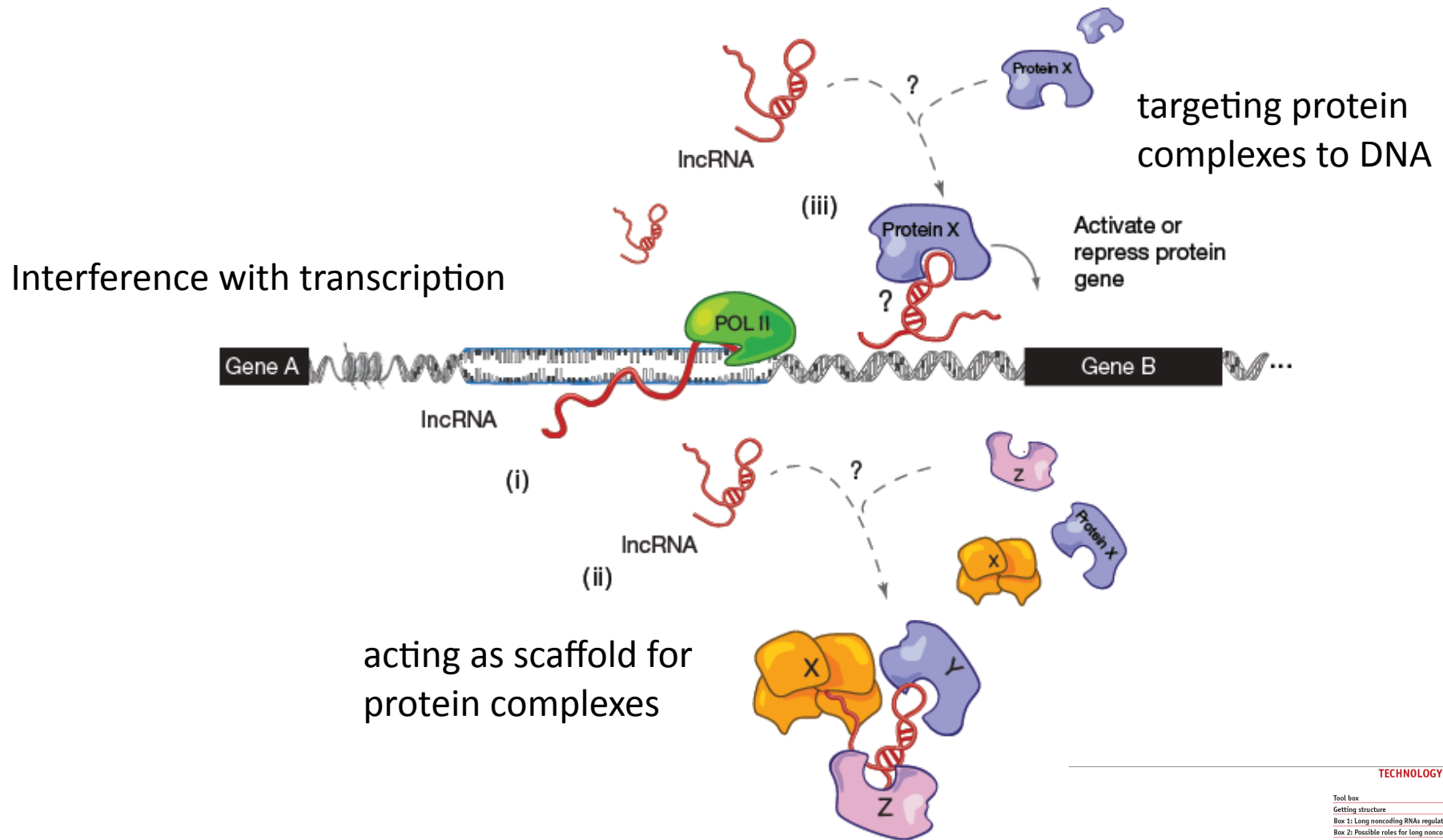


Post-transcriptional processing



Long non-coding RNAs:
insights into functions

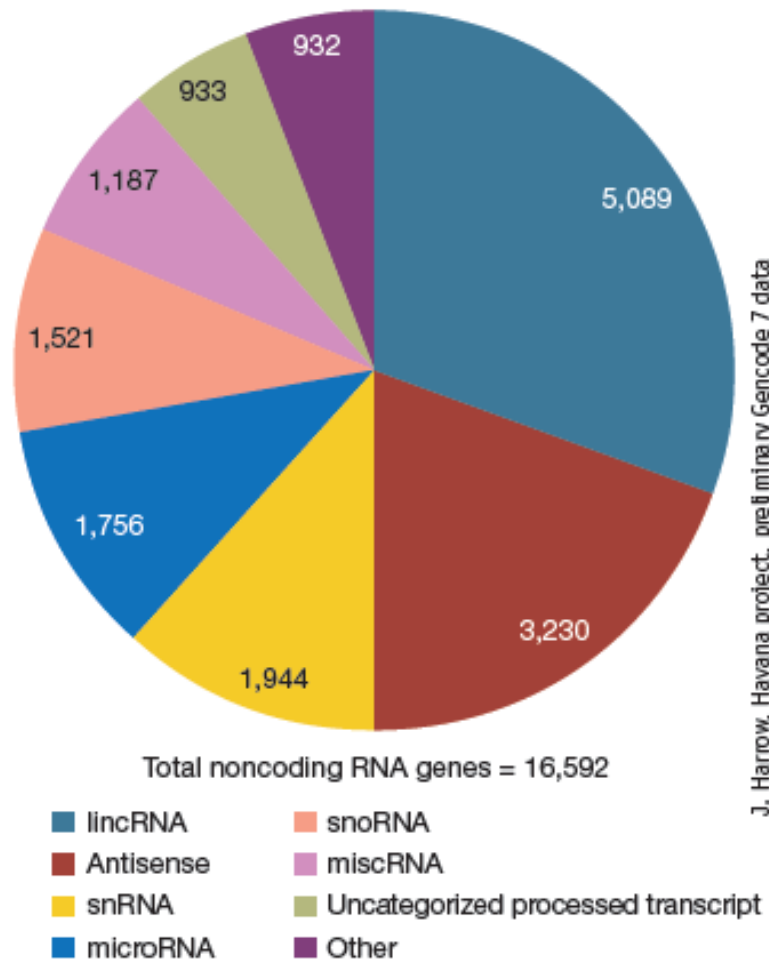
Molecular mechanisms of how ncRNAs can function



TECHNOLOGY FEATURE

Tool box	380
Getting structure	382
Box 1: Long noncoding RNAs regulating genes	380
Box 2: Possible roles for long noncoding RNA	381

Distribution of ncRNAs



long ncRNAs are similar to mRNAs:

- RNA polymerase II-promoted
- polyadenylated
- often alternatively spliced

TECHNOLOGY FEATURE

Tool box	380
Getting structure	382
Box 1: Long noncoding RNAs regulating genes	380
Box 2: Possible roles for long noncoding RNA	381

Long noncoding RNAs: the search for function

Monya Baker

FOCUS OF THIS REVIEW

Focus of this review

Many studies have identified novel ncRNAs, tried to classify them, and assign biological functions to them. Various approaches have been developed these steps.

Here I will review the different approaches to

1. identify ncRNAs
2. assign a biological function or molecular mechanism to ncRNAs

To get a better idea of how different studies can be compared

Very general overview of approaches to identify ncRNAs

Annotating non-coding transcription using functional genomics strategies

Alistair R. R. Forrest, Rehab F. Abdelhamid and Piero Carninci

Abstract

Non-coding RNA (ncRNA) transcripts are RNA molecules that do not code for proteins, but elicit function by other mechanisms. The vast majority of RNA produced in a cell is non-coding ribosomal RNA, produced from relatively few loci, however more recently complementary DNA (cDNA) cloning, tag sequencing, and genome tiling array studies suggest that ncRNAs also account for the majority of RNA species produced by a cell. ncRNA based regulation has been referred to as a 'hidden layer' of signals or 'dark matter' that control gene expression in cellular processes by poorly described mechanisms. These terms have appeared as ncRNAs until recently have been ignored by expression profiling and cDNA annotation projects and their mode of action is diverse (e.g. influencing chromatin structure and epigenetics, translational silencing, transcriptional silencing). Here, we highlight recent functional genomics strategies toward identifying and assigning function to ncRNA transcription.

Keywords: non-coding RNA; Sequencing; transcription; annotation

IDENTIFICATION OF NON-CODING RNA

1. Identification and validation of functional ncRNA

a) Identification

- I. Identification based on template-free expression measurement
- II. Identification based on chromatin-state signatures
- III. Identification of unstable transcripts

b) Validation and how to distinguish coding from non-coding

- I. conservation
- II. distinguish coding from non-coding RNA

c) Application: Computational identification based on published data

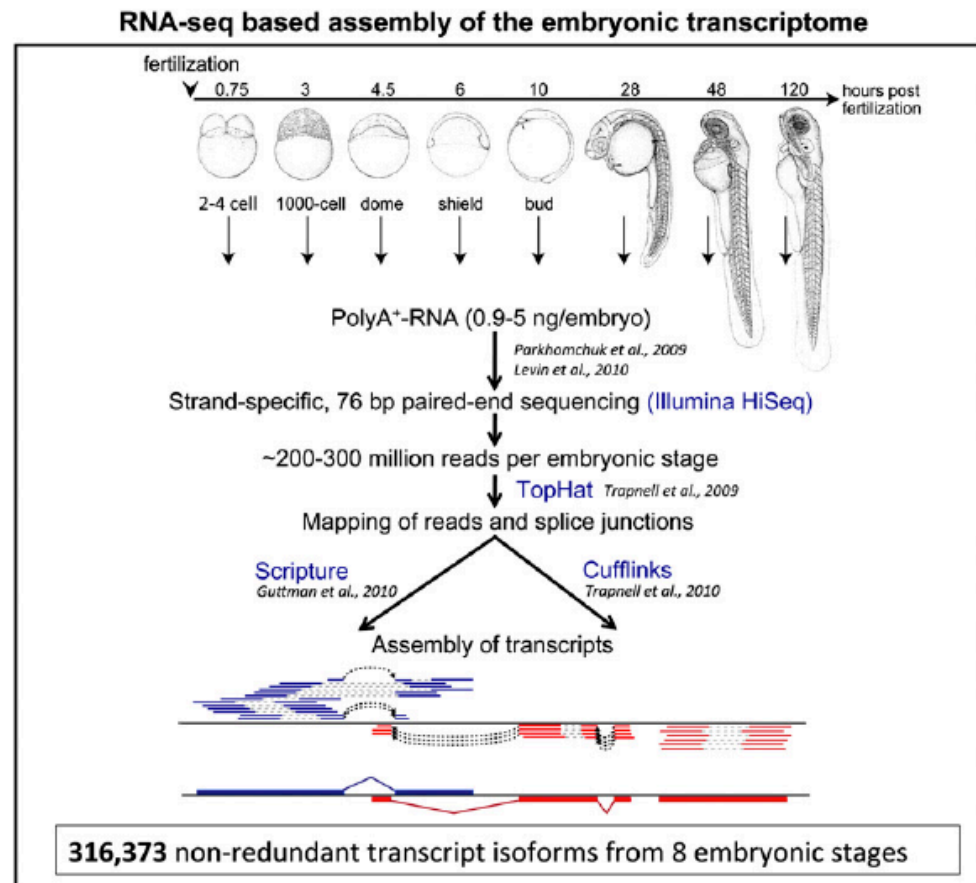
- RNA-seq
- tiling arrays
- global run-on (GRO)-seq: detection of nascent transcripts

a) I. RNA-seq: zebrafish

Resource

Systematic identification of long noncoding RNAs expressed during zebrafish embryogenesis

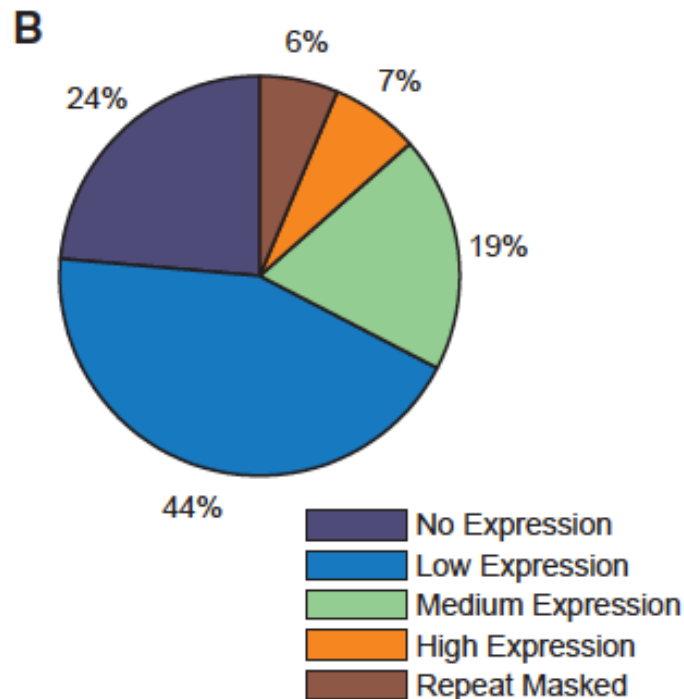
Andrea Pauli,^{1,7,8} Eivind Valen,^{2,7} Michael F. Lin,^{3,4} Manuel C Nadine L. Vastenhouw,¹ Joshua Z. Levin,⁴ Lin Fan,⁴ Albin Sar John L. Rinn,^{4,5} Aviv Regev,^{3,4,6,8} and Alexander F. Schier^{1,4,8}



a) I. RNA-seq: first studies in yeast

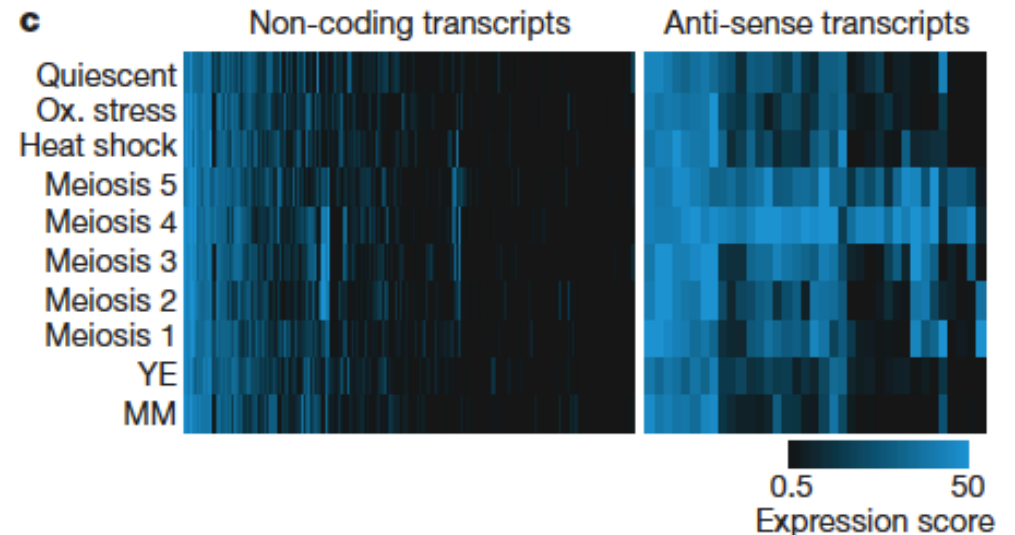
The Transcriptional Landscape of the Yeast Genome Defined by RNA Sequencing

Ugrappa Nagalakshmi,^{1*} Zhong Wang,^{1*} Karl Waern,¹ Chong Shou,² Debasish Raha,¹ Mark Gerstein,^{2,3} Michael Snyder^{1,2,3†}



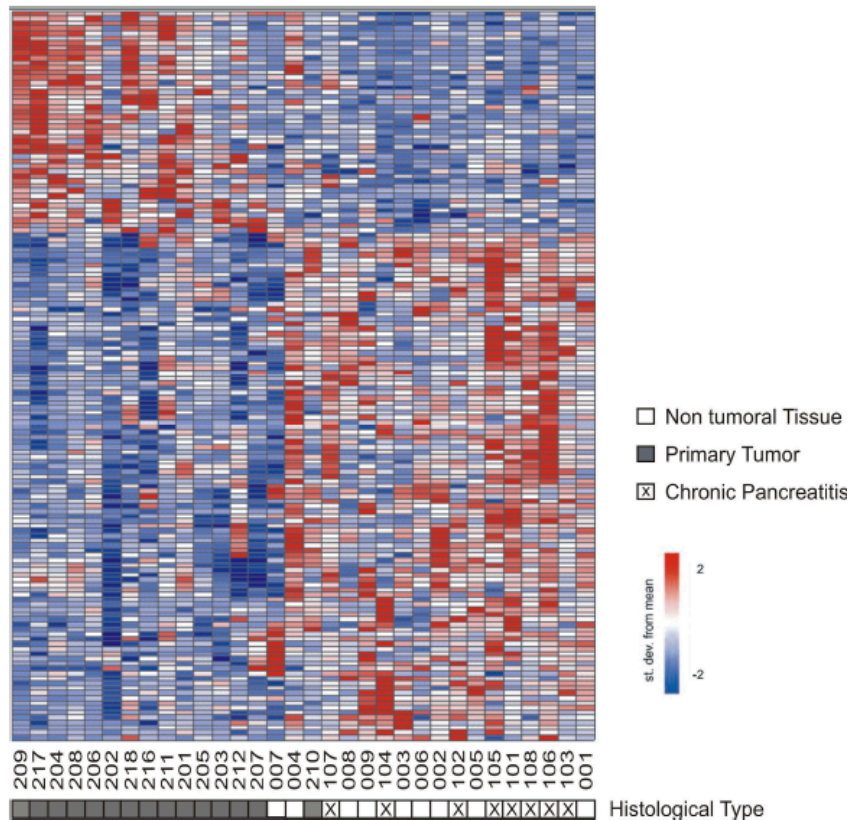
Dynamic repertoire of a eukaryotic transcriptome surveyed at single-nucleotide resolution

Brian T. Wilhelm^{1*†}, Samuel Marguerat^{1*†}, Stephen Watt^{1†}, Falk Schubert^{1†}, Valerie Wood¹, Ian Goodhead^{1†}, Christopher J. Penkett^{1†}, Jane Rogers¹ & Jürg Bähler^{1†}



Long noncoding intronic RNAs are differentially expressed in primary and metastatic pancreatic cancer

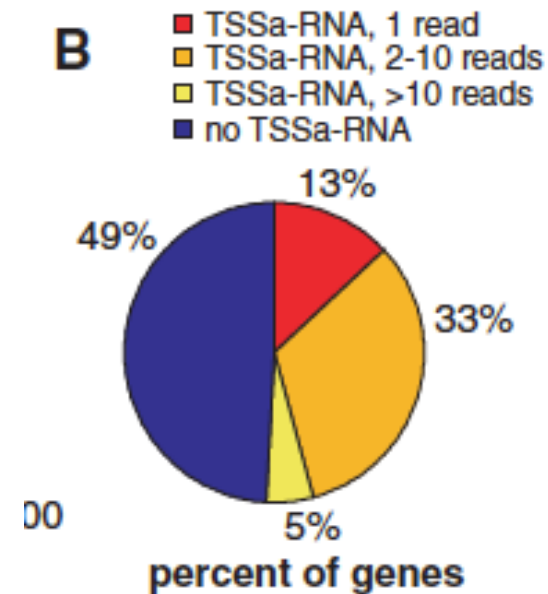
Ana C Tahira¹, Márcia S Kubrusly², Michele F Faria¹, Bianca Dazzani¹, Rogério S Fonseca¹, Vinicius Maracaja-Coutinho¹, Sergio Verjovski-Almeida¹, Marcel CC Machado² and Eduardo M Reis^{1*}



a) I. RNA-seq: humans

Divergent Transcription from Active Promoters

Amy C. Seila,^{1*} J. Mauro Calabrese,^{1,2*}† Stuart S. Levine,³ Gene W. Yeo,⁴‡ Peter B. Rahl,³ Ryan A. Flynn,¹ Richard A. Young,^{2,3} Phillip A. Sharp^{1,2,5}



RNA-seq for identifying ncRNA

+

- identification of novel transcripts
- in principle reference-independent

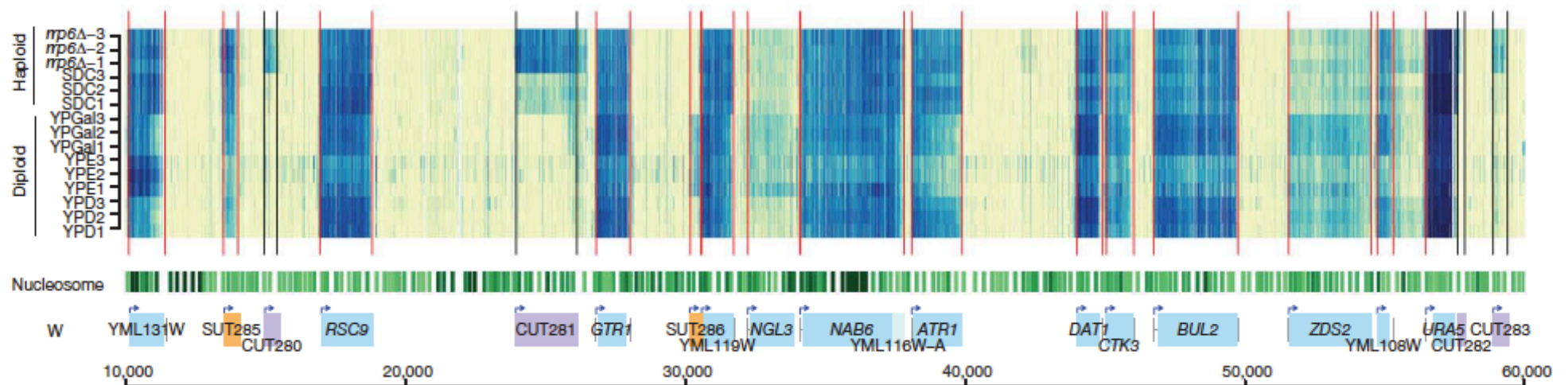
-

- no detection of unstable transcripts
- need high coverage for detecting ncRNAs (they are generally lowly expressed)

a) I. tiling-array

Bidirectional promoters generate pervasive transcription in yeast

Zhenyu Xu^{1*}, Wu Wei^{1*}, Julien Gagneur¹, Fabiana Perocchi¹, Sandra Clauder-Münster¹, Jurgi Camblong²,



manual annotation of novel transcripts

tilling-arrays for identifying ncRNA

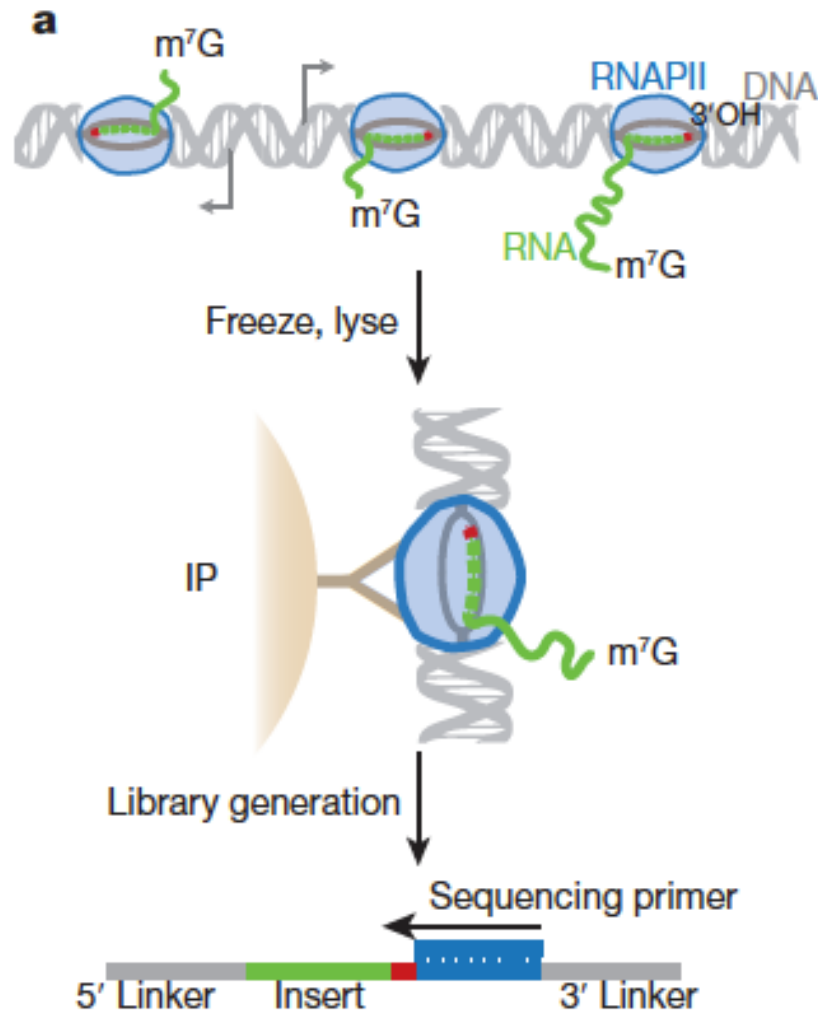
+

- identification of novel transcripts
- for whole genome arrays: reference-independent

-

- no detection of unstable transcripts
- potential difficulty to detect low abundance transcripts (they might fall into the background region)

a) I. GRO-seq (Global run-on)



Nascent transcript sequencing visualizes transcription at nucleotide resolution

L. Stirling Churchman¹ & Jonathan S. Weissman¹

GRO-seq for identifying ncRNA

+

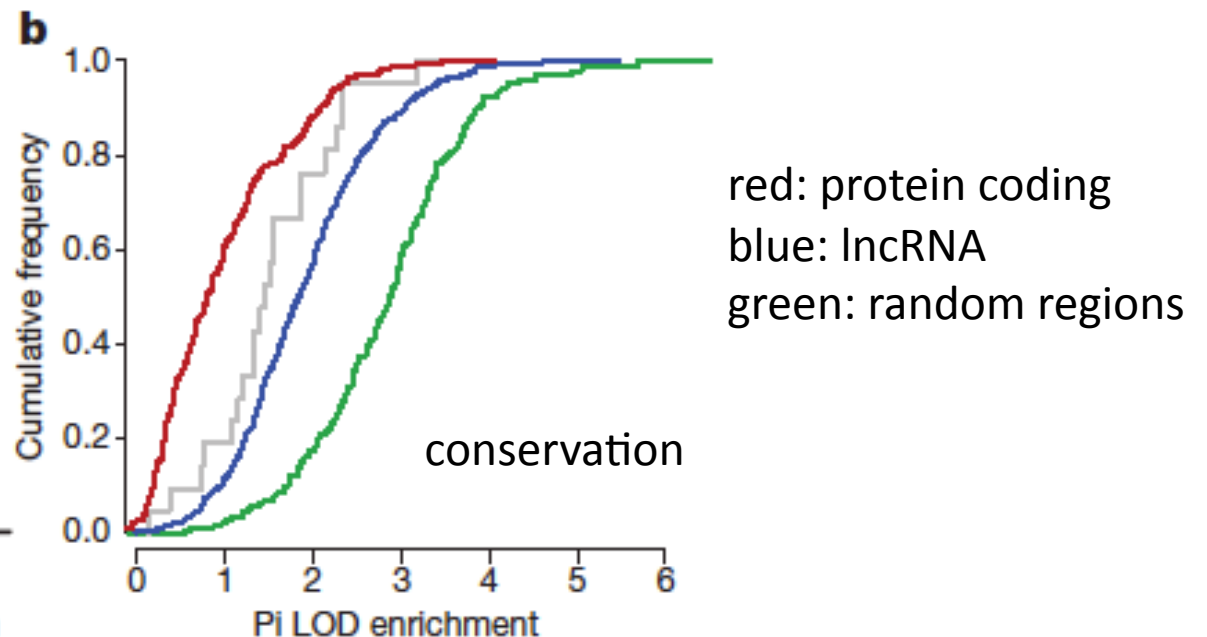
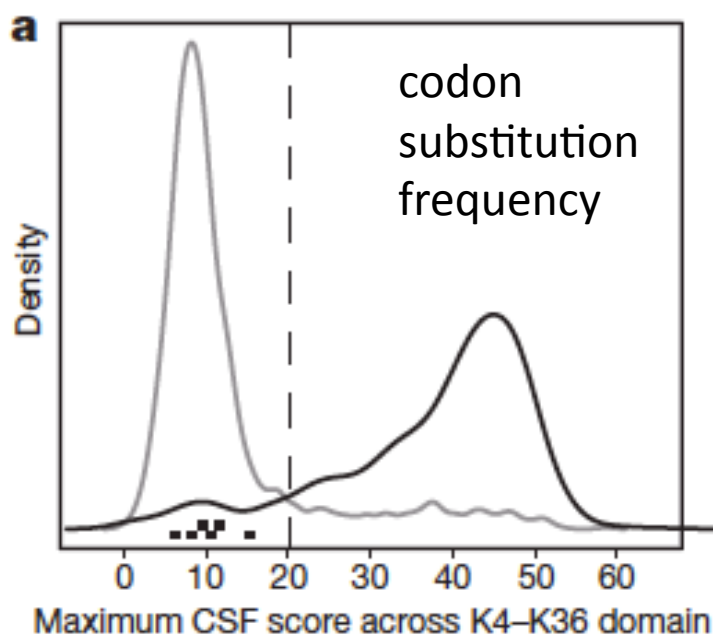
- identification of novel transcripts
- reference-independent
- detection of unstable transcripts

-

- difficult to obtain the transcript lengths
- only unspliced RNAs detected (except for when co-transcriptional splicing occurs)

a) II. Identification based on chromatin state signatures

- active transcription: H3K4me3 at promoter, H3K36me3 across body
- Find K4–K36 regions outside coding genes loci to discover ncRNAs



Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals

Mitchell Guttman^{1,2}, Ido Amit¹, Manuel Garber¹, Courtney French¹, Michael F. Lin¹, David Feldser³, Maite Huarte^{1,6}, Or Zuk¹, Bryce W. Carey^{2,8}, John P. Cassidy^{2,8}, Moran N. Cabili⁷, Rudolf Jaenisch^{2,8}, Tarjei S. Mikkelsen^{1,4}, Tyler Jacks^{2,3}, Nir Hacohen^{1,5}, Bradley E. Bernstein^{1,10,11}, Manolis Kellis^{1,5}, Aviv Regev^{1,2}, John L. Rinn^{1,6,11*} & Eric S. Lander^{1,2,7,8*}

Chromatin ChIP for identifying ncRNA

+

- identification of regions that are transcribed
- reference-independent (if combined with ChIP-seq)
- detection of unstable transcripts if they have same chromatin state

-

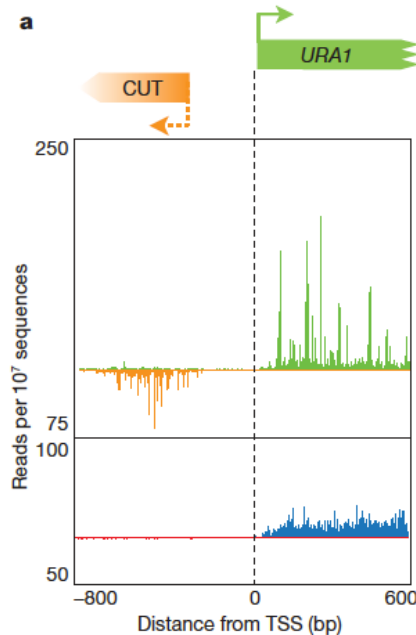
- indirect measure of expression

a) III: identification of unstable transcripts

- GRO-seq
- exosome knockout/
knock-down

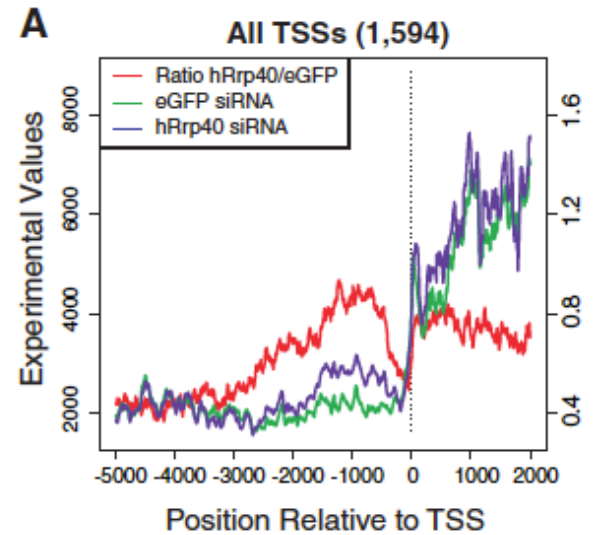
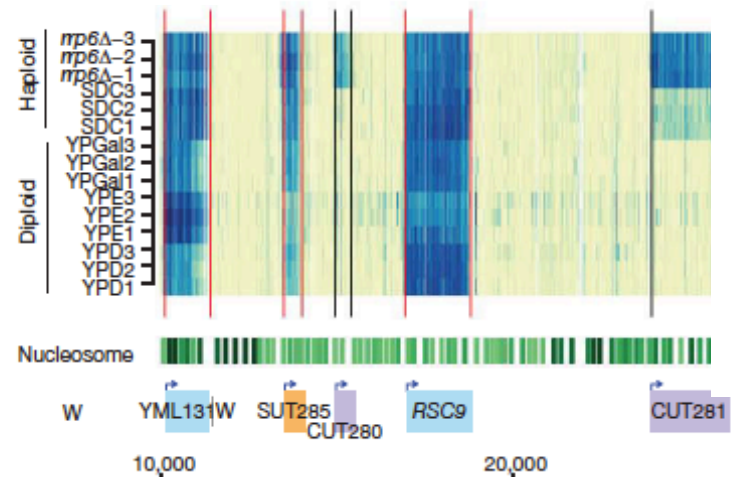
Nascent transcript sequencing visualizes transcription at nucleotide resolution

L. Stirling Churchman¹ & Jonathan S. Weissman¹



RNA Exosome Depletion Reveals Transcription Upstream of Active Human Promoters

Pascal Preker,¹ Jesper Nielsen,² Susanne Kammler,^{1*} Søren Lykke-Andersen,¹ Marianne S. Christensen,¹ Christoph



Bidirectional promoters generate pervasive transcription in yeast

Zhenyu Xu^{1*}, Wu Wei^{1*}, Julien Gagneur¹, Fabiana Perocchi¹, Sandra Clauder-Münster¹, Jurgi Camblong², Elisa Guffanti³, Françoise Stutz³, Wolfgang Huber⁴ & Lars M. Steinmetz¹

identification of unstable transcripts

+/- points for the individual techniques apply in addition

+

- comprehensive identification of transcripts

-

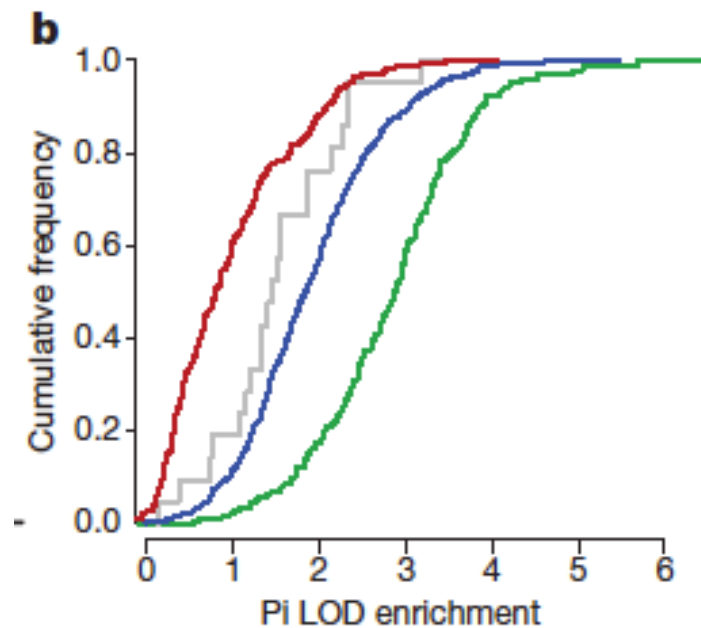
- challenge to assign functional significance

b) validation of non-coding transcripts

- I. Are they functional
- II. Are they non-coding

b) I. Validation: is the ncRNA functional

conservation: ncRNA are more conserved than random regions



red: protein coding
blue: lncRNA
green: random regions

Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals

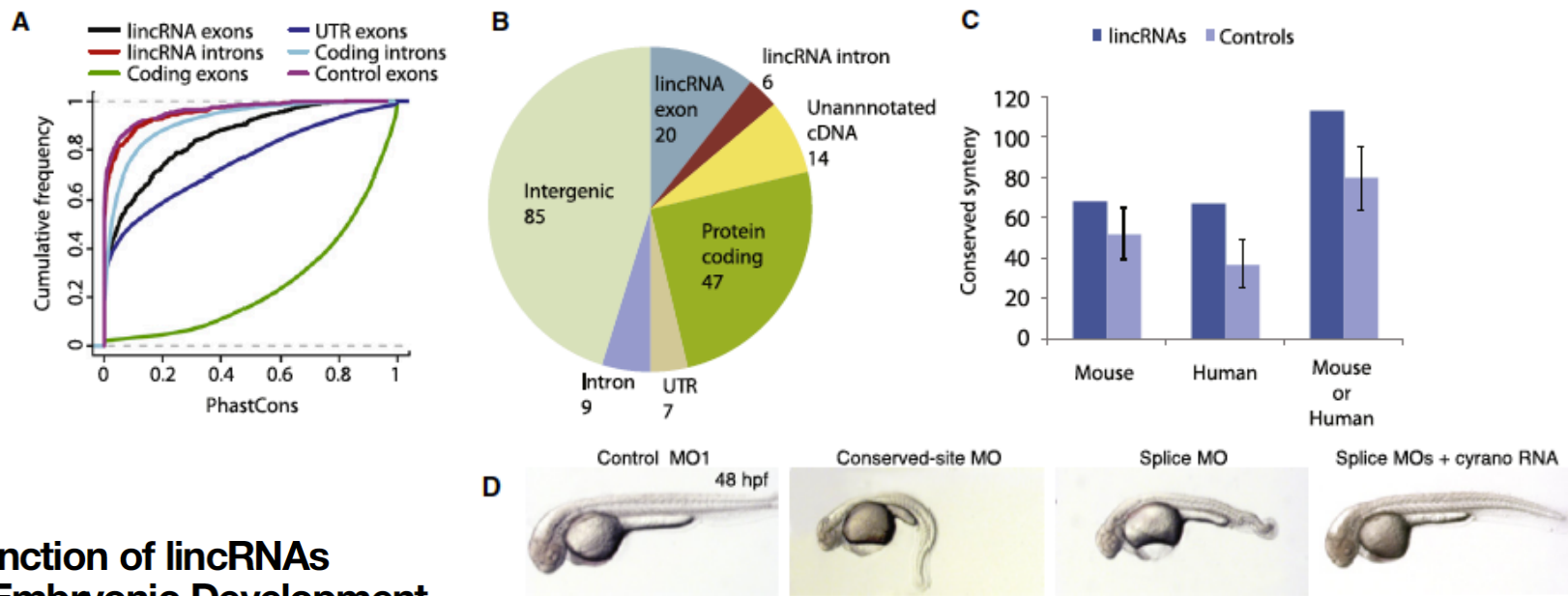
Mitchell Guttman^{1,2}, Ido Amit¹, Manuel Garber¹, Courtney French¹, Michael F. Lin¹, David Feldser³, Maite Huarte^{1,6}, Or Zuk¹, Bryce W. Carey^{2,8}, John P. Cassidy^{2,8}, Moran N. Cabili⁷, Rudolf Jaenisch^{2,8}, Tarjei S. Mikkelsen^{1,4}, Tyler Jacks^{2,9}, Nir Hacohen¹⁰, Bradley E. Bernstein^{10,11}, Manolis Kellis¹², Aviv Regev¹², John L. Rinn^{1,6,11} & Eric S. Lander^{1,2,7,8*}

b) I. Validation: is the ncRNA functional - caveat

ncRNA more conserved than random regions

BUT location more conserved than sequence:

—e.g. possible to rescue ncRNA knockout zebrafish with mouse ncRNA

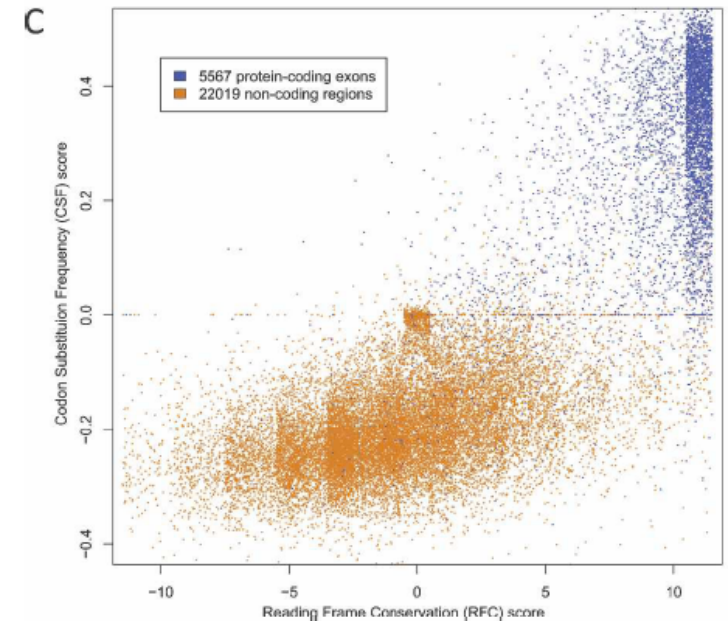
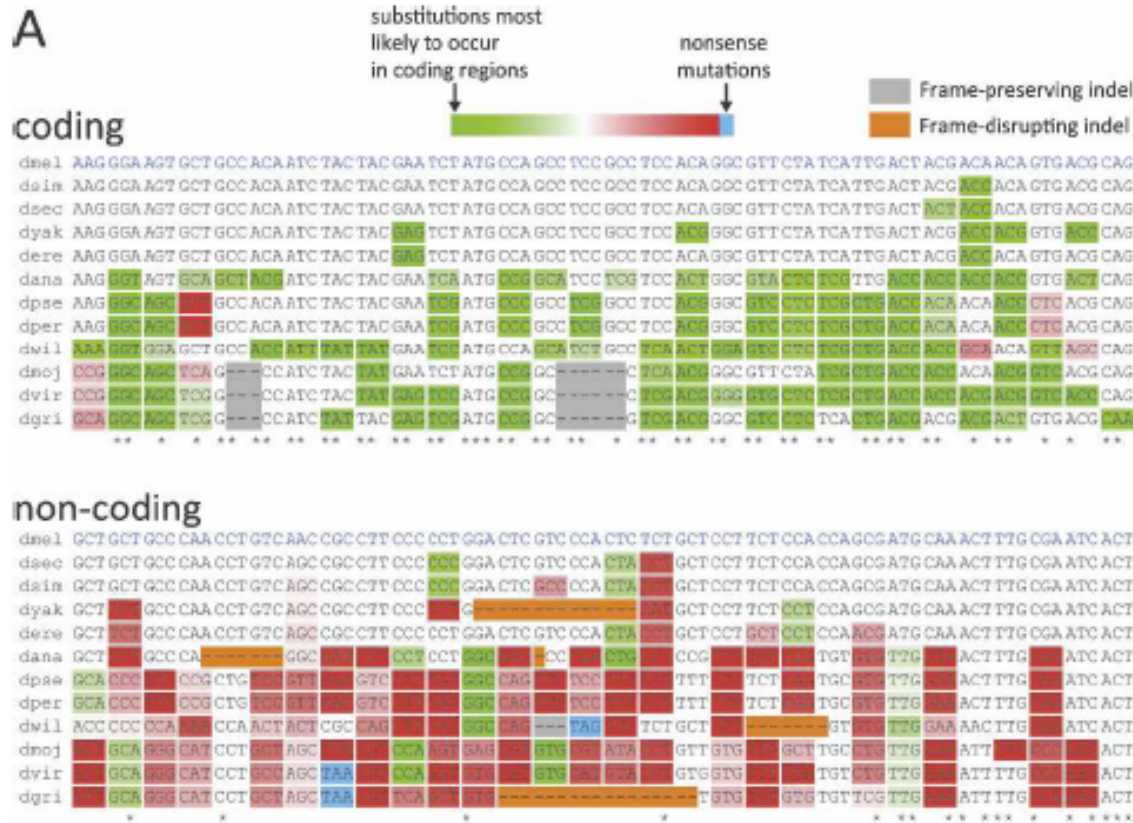


Conserved Function of lincRNAs in Vertebrate Embryonic Development despite Rapid Sequence Evolution

b) Validation: II. is the transcript non-coding?

- BLASTP the cDNAs sequences
- Assessment of codon substitution frequency
- Assessment of reading frame conservation

codon substitution frequency and reading frame conservation



12 *Drosophila* Genomes/Letter

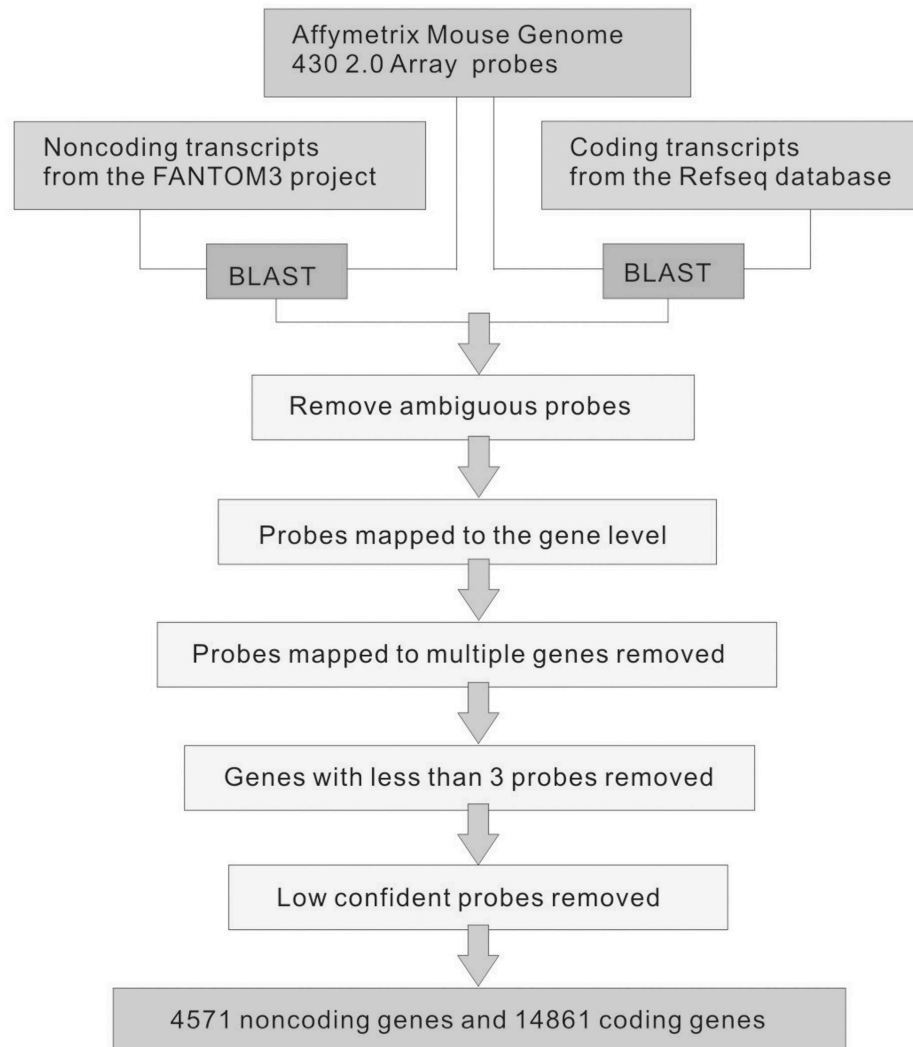
Revisiting the protein-coding gene catalog of *Drosophila melanogaster* using 12 fly genomes

Michael F. Lin,¹ Joseph W. Carlson,² Madeline A. Crosby,³ Beverley B. Matthews,³ Charles Yu,² Soo Park,² Kenneth H. Wan,² Andrew J. Schroeder,³ L. Sian Gramates,³ Susan E. St. Pierre,³ Margaret Roark,³ Kenneth L. Wiley Jr.,⁴ Rob J. Kulathinal,³ Peili Zhang,³ Kyle V. Myrick,⁴ Jerry V. Antone,⁴ Susan E. Celniker,² William M. Gelbart,^{3,4} and Manolis Kellis^{1,3,6}

- Within coding regions, triplet substitutions are biased toward conservative codon substitutions (Codon Substitution Frequencies, CSF).
- Indels in coding regions are strongly biased to be a multiple of three in length (reading frame conservation; RFC).

c) Computational identification based on published data: using microarray data – remapping the probes and literature curation

A



BIOINFORMATICS

Genome-wide computational identification and manual annotation of human long noncoding RNA genes

HUI JIA,¹ MAUREEN OSAK,² GIREESH K. BOGU,³ LAWRENCE W. STANTON,³ RORY JOHNSON,^{3,4} and LEONARD LIPOVICH¹

¹Center for Molecular Medicine and Genetics, Wayne State University, Detroit, Michigan 48202, USA

²Lee and Roland Witte Natural Sciences Division, Hillsdale College, Hillsdale, Michigan 49242, USA

³Stem Cell and Developmental Biology Group, Genome Institute of Singapore, 138672 Singapore

D210–D215 *Nucleic Acids Research*, 2012, Vol. 40, Database issue
doi:10.1093/nar/gkr1175

Published online 1 December 2011

NONCODE v3.0: integrative annotation of long noncoding RNAs

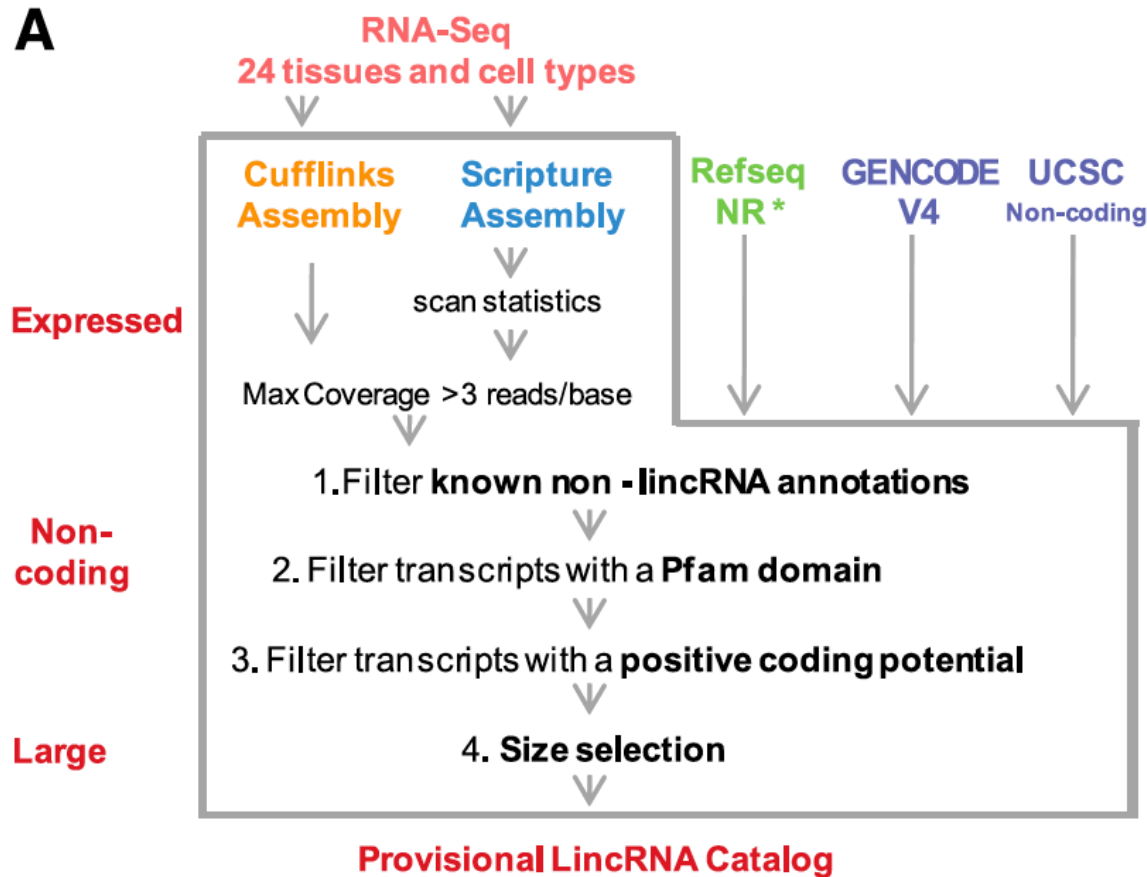
Dechao Bu^{1,2}, Kuntao Yu^{1,2}, Silong Sun¹, Chaoyong Xie^{1,2}, Geir Skogerboe³, Ruoyu Miao^{1,4}, Hui Xiao¹, Qi Liao¹, Haitao Luo¹, Guoguang Zhao^{1,2}, Haitao Zhao⁴, Zhiyong Liu¹, Changning Liu¹, Runsheng Chen^{3,*} and Yi Zhao^{1,*}

W118–W124 *Nucleic Acids Research*, 2011, Vol. 39, Web Server issue
doi:10.1093/nar/gkr432

ncFANs: a web server for functional annotation of long non-coding RNAs

Qi Liao^{1,2,3}, Hui Xiao¹, Dechao Bu^{1,4}, Chaoyong Xie¹, Ruoyu Miao⁵, Haitao Luo¹, Guoguang Zhao^{1,4}, Kuntao Yu^{1,4}, Haitao Zhao⁵, Geir Skogerboe⁶, Runsheng Chen⁶, Zhongdao Wu^{2,3}, Changning Liu^{1,*} and Yi Zhao^{1,*}

c) Computational identification based on published data: using RNA-seq data



Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses

Moran N. Cabili,^{1,2,3} Cole Trapnell,^{1,3} Loyal Goff,^{1,4} Magdalena Koziol,^{1,3} Barbara Tazon-Vega,^{1,3} Aviv Regev,^{1,3,6} and John L. Rinn^{1,3,6,7}

Summary ncRNA identification

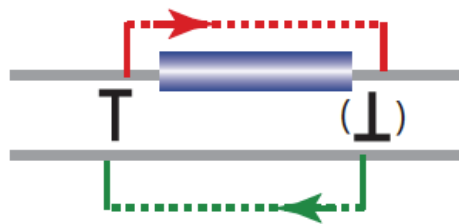
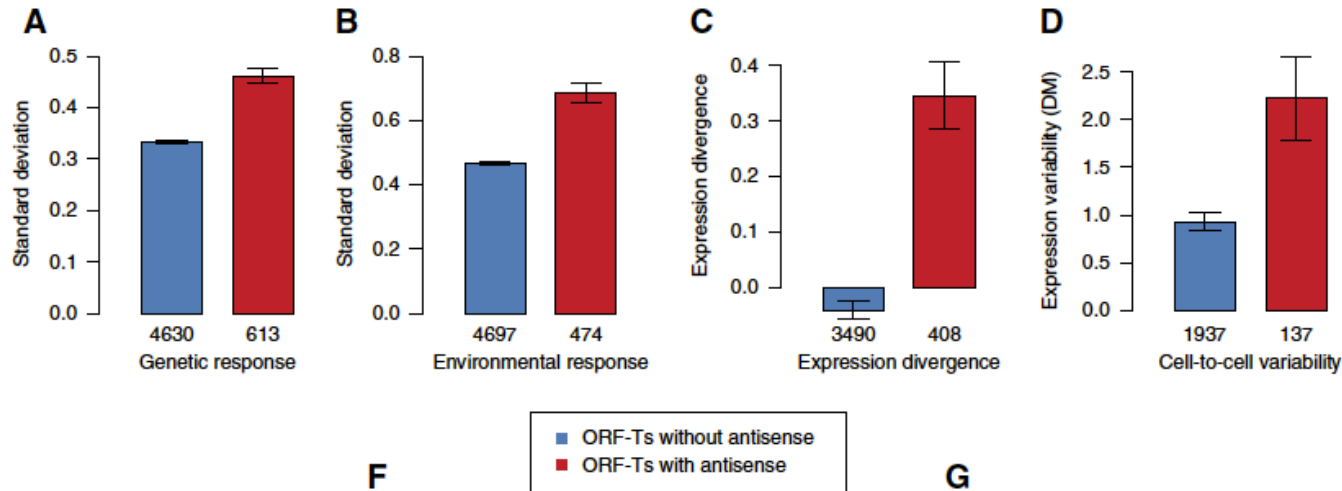
- Different template-free methods exist for detecting non-coding RNAs
- Methods to validate ncRNAs include
 - conservation: less than proteins, more than random – use positional conservation
 - codon substitution frequency
 - reading frame conservation
- Applications: meta studies integrating various published data sets to identify novel ncRNAs

FUNCTIONAL IMPACT

Inferring the functional role of ncRNAs

- a) Molecular mechanisms of ncRNAs activity
 - I. sense antisense regulation
 - II. association with protein complexes
 - III. perturbation studies
- b) Biological function
 - I. tissue specificity
 - II. functional analysis of genes co-expressed with ncRNAs

a) I. Sense antisense regulation: yeast



Molecular Systems Biology 7; Article number 468; doi:10.1038/msb.2011.1
 Citation: *Molecular Systems Biology* 7:468
 © 2011 EMBO and Macmillan Publishers Limited. All rights reserved 1744-4292/11
 www.molecular-systems-biology.com

molecular
systems
biology

Antisense expression increases gene expression variability and locus interdependency

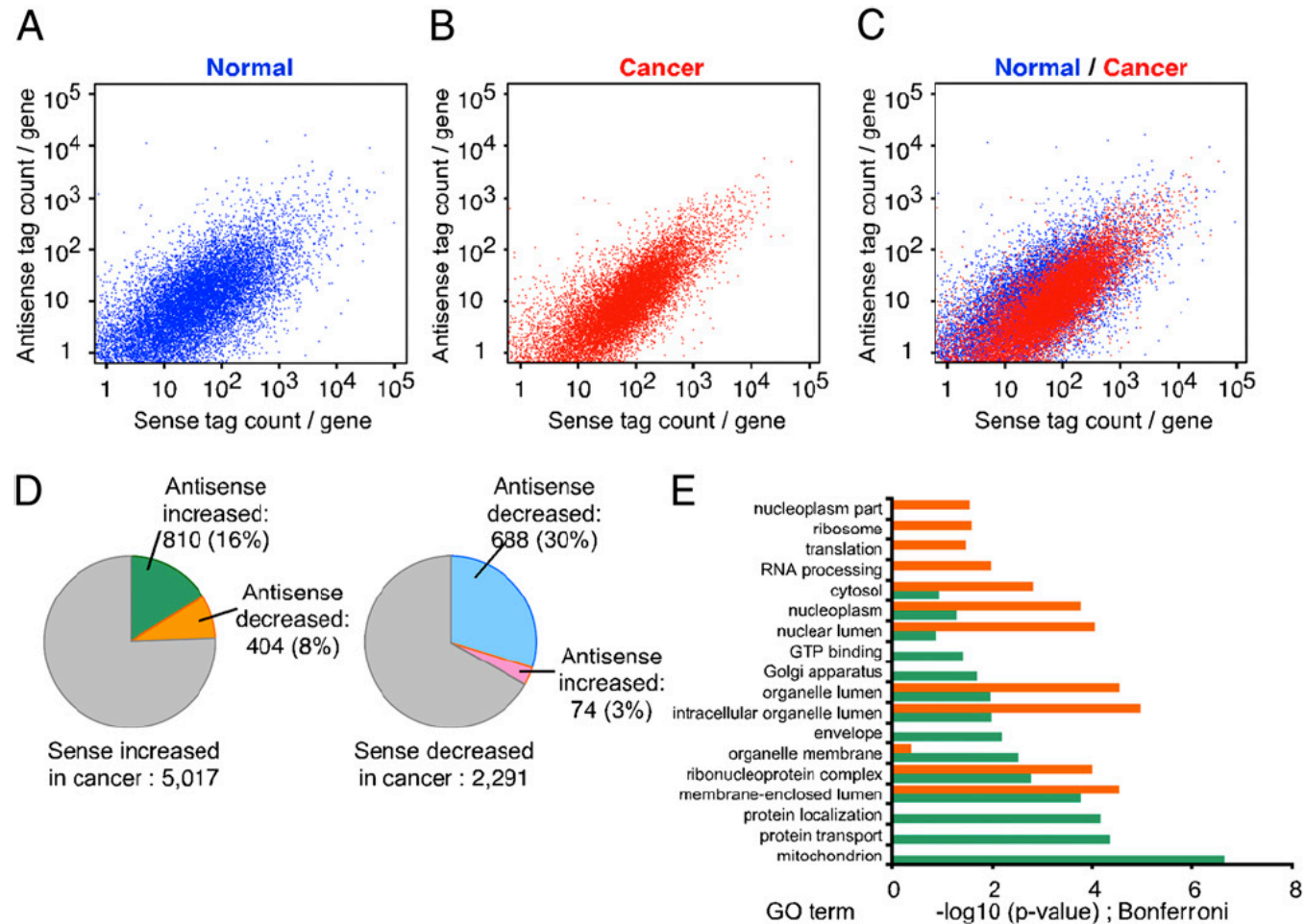
Zhenyu Xu¹, Wu Wei¹, Julien Gagneur¹, Sandra Clauder-Münster, Miłosz Smolik, Wolfgang Huber and Lars M Steinmetz*

Genome Biology Unit, European Molecular Biology Laboratory, Heidelberg, Germany

¹ These authors contributed equally to this work

* Corresponding author. Genome Biology Unit, European Molecular Biology Laboratory, Meyerhofstrasse 1, Heidelberg 69117, Germany.
 Tel: +49 6221 387 8389; Fax: +49 6221 387 8518; E-mail: larsms@embl.de

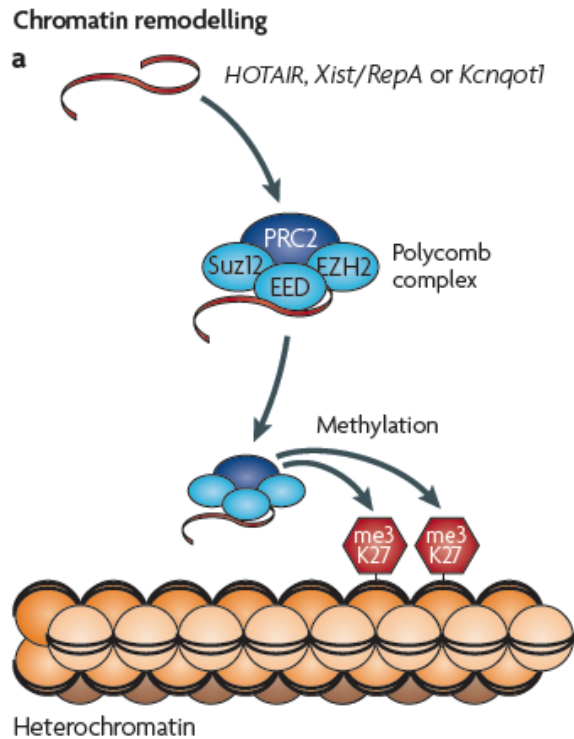
a) I. Sense antisense regulation: breast cancer



Altered antisense-to-sense transcript ratios in breast cancer

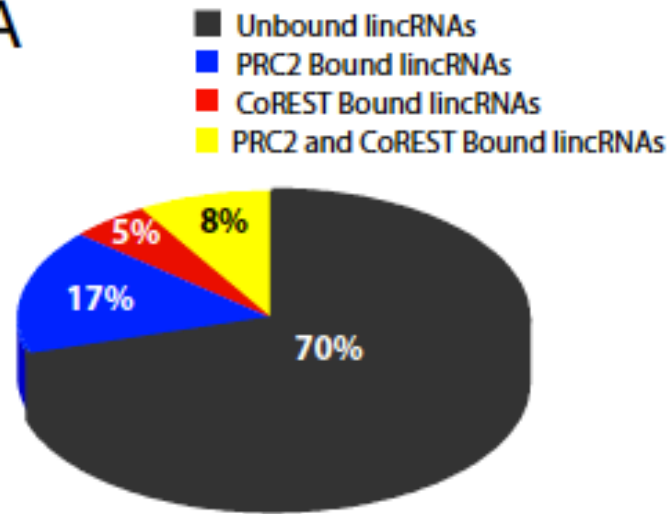
Reo Maruyama^{a,b,c,1}, Michail Shiptsin^{a,b,c,1}, Sibgat Choudhury^{a,b,c,1}, Zhenhua Wu^{d,e,1}, Alexei Protopopov^f, Jun Yao^g, Pang-Kuo Lo^h, Marina Bessarabovaⁱ, Alex Ishkinⁱ, Yuri Nikolsky^j, X. Shirley Liu^{d,e}, Saraswati Sukumar^h, and Kornelia Polyak^{a,b,c,k,2}

a) II. association with protein complexes



Inspired by example

A



Protein-coding genes associated with PRC2
<2%



Many human large intergenic noncoding RNAs associate with chromatin-modifying complexes and affect gene expression

Ahmad M. Khalil^{a,b,1}, Mitchell Guttman^{a,c,1}, Maite Huarte^{a,b}, Manuel Garber^a, Arjun Raj^d, Dianali Rivea Morales^{a,b}, Celly Thomas^{a,b}, Aviva Presser^a, Bradley E. Bernstein^{a,e}, Alexander van Oudenaarden^d, Aviv Regev^{a,c}, Eric S. Lander^{a,c,f,1,2}, and John L. Rinn^{a,b,1,2}

The Broad Institute of Harvard and Massachusetts Institute of Technology, Cambridge, MA 02142; ¹Department of Pathology, Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, MA 02115; ²Department of Biology and ³Genetics, Massachusetts Institute of Technology, Cambridge, MA 02139

a) III. perturbation studies: special case yeast

- ncRNAs exported to cytoplasm
- ncRNAs specifically stabilized in meiosis required
- ncRNAs only detectable in exosome knockouts

LETTER

doi:10.1038/nature10118

XUTs are a class of Xrn1-sensitive antisense regulatory non-coding RNA in yeast

E. L. van Dijk^{1*}, C. L. Chen^{1*}, Y. d'Aubenton-Carafa¹, S. Gourvennec², M. Kwapisz², V. Roche², C. Bertrand², M. Silvain¹, P. Legoux-Né³, S. Loillet⁴, A. Nicolas⁴, C. Thermes¹ & A. Morillon^{1,2}

Execution of the meiotic noncoding RNA expression program and the onset of gametogenesis in yeast require the conserved exosome subunit Rrp6

Aurélie Lardenois^{a,1}, Yuchen Liu^{a,1}, Thomas Walther^b, Frédéric Chalmel^a, Bertrand Evrard^a, Marina Granovskaia^c, Angela Chu^d, Ronald W. Davis^{d,e,2}, Lars M. Steinmetz^c, and Michael Primig^{a,2}

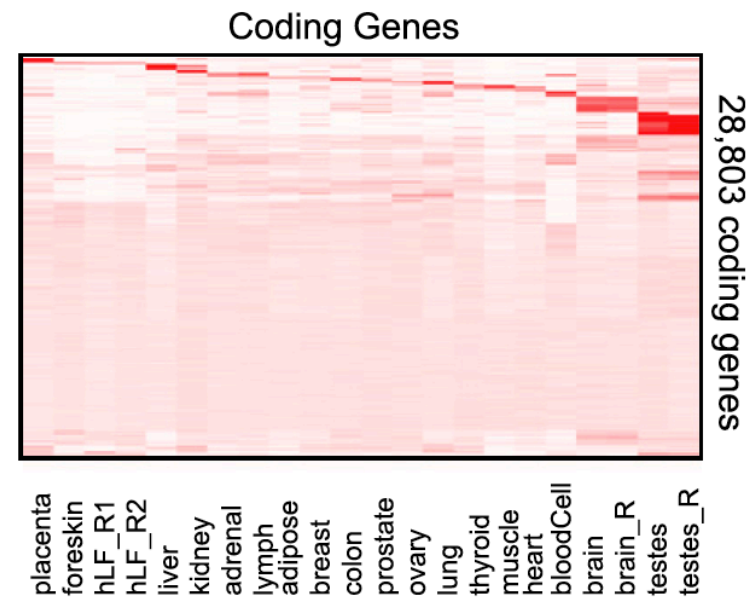
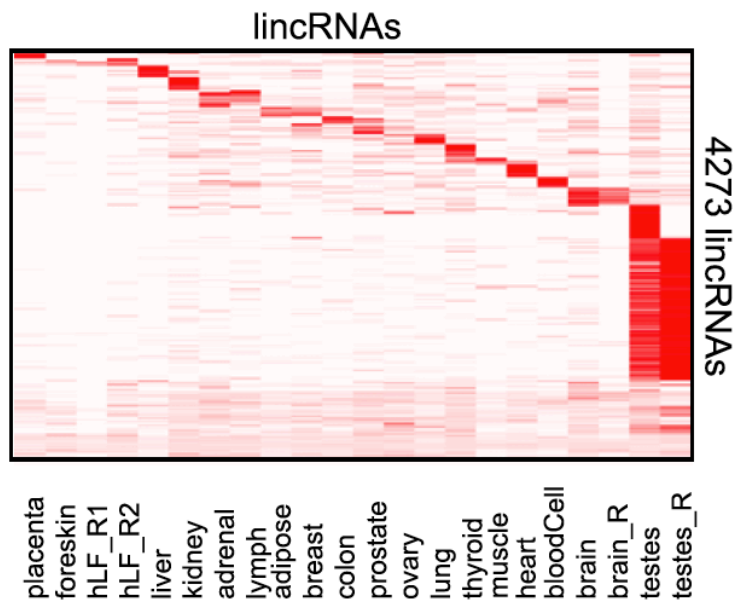
LETTERS

Bidirectional promoters generate pervasive transcription in yeast

Zhenyu Xu^{1*}, Wu Wei^{1*}, Julien Gagneur¹, Fabiana Perocchi¹, Sandra Clauder-Münster¹, Jurgi Camblong², Elisa Guffanti³, Françoise Stutz³, Wolfgang Huber⁴ & Lars M. Steinmetz¹

b) I. Tissue-specificity

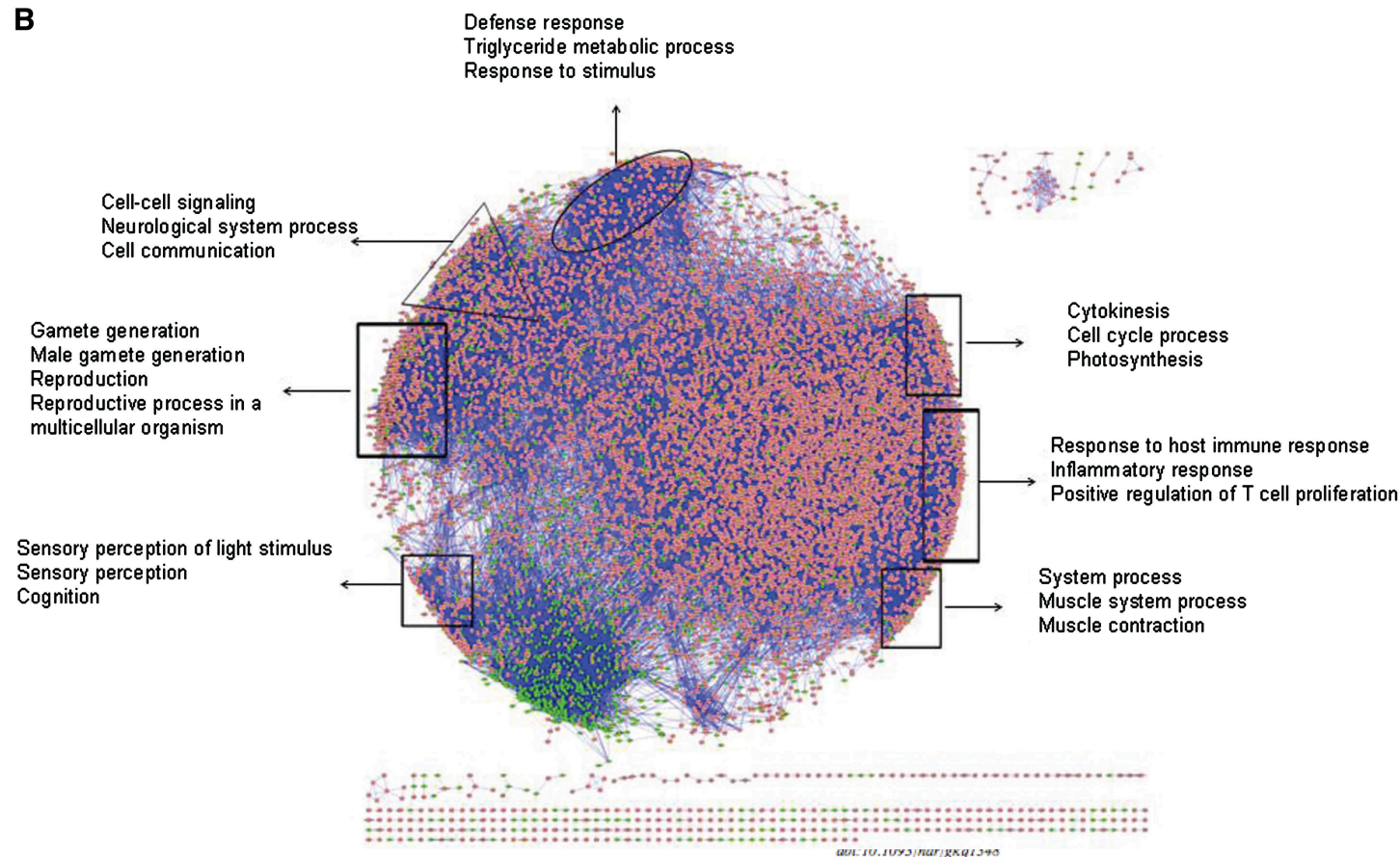
A



Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses

Moran N. Cabili,^{1,2,3} Cole Trapnell,^{1,3} Loyal Goff,^{1,4} Magdalena Koziol,^{1,3} Barbara Tazon-Vega,^{1,3} Aviv Regev,^{1,3,6} and John L. Rinn^{1,3,6,7}

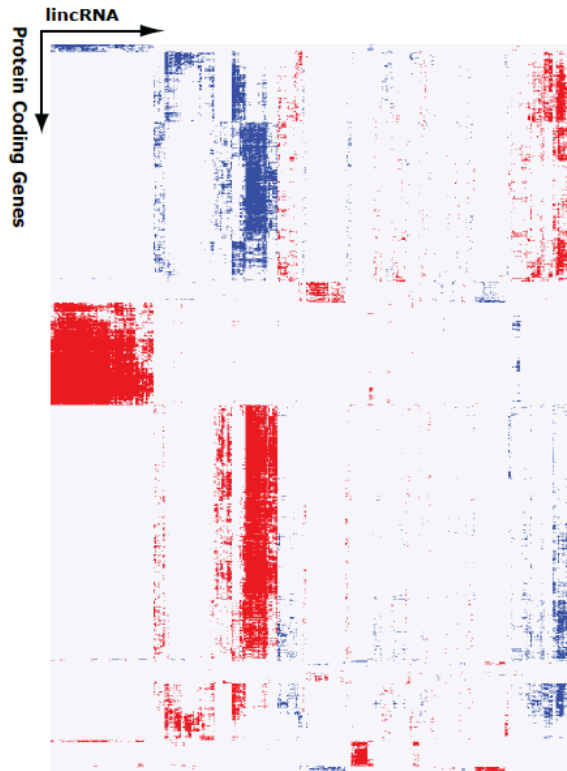
b) II. functional analysis of genes co-expressed with ncRNAs: pairwise correlation with genes



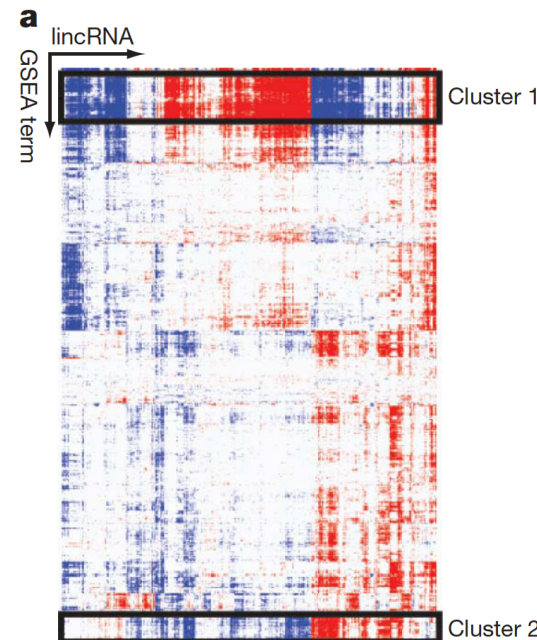
Large-scale prediction of long non-coding RNA functions in a coding–non-coding gene co-expression network

Qi Liao^{1,2,3}, Changning Liu¹, Xiongying Yuan^{1,4}, Shuli Kang¹, Ruoyu Miao⁵, Hui Xiao¹, Guoguang Zhao^{1,4}, Haitao Luo¹, Dechao Bu^{1,4}, Haitao Zhao⁵, Geir Skogerboe⁶, Zhongdao Wu^{2,3,*} and Yi Zhao^{1,*}

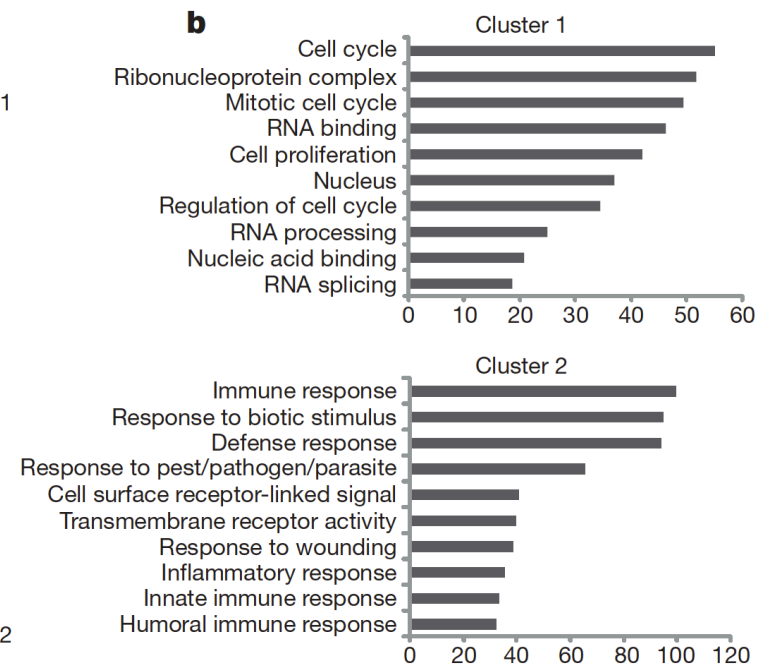
Correlation of ncRNA and proteins across 19 conditions



Correlation ncRNA with gene sets (sets of proteins) and clustering



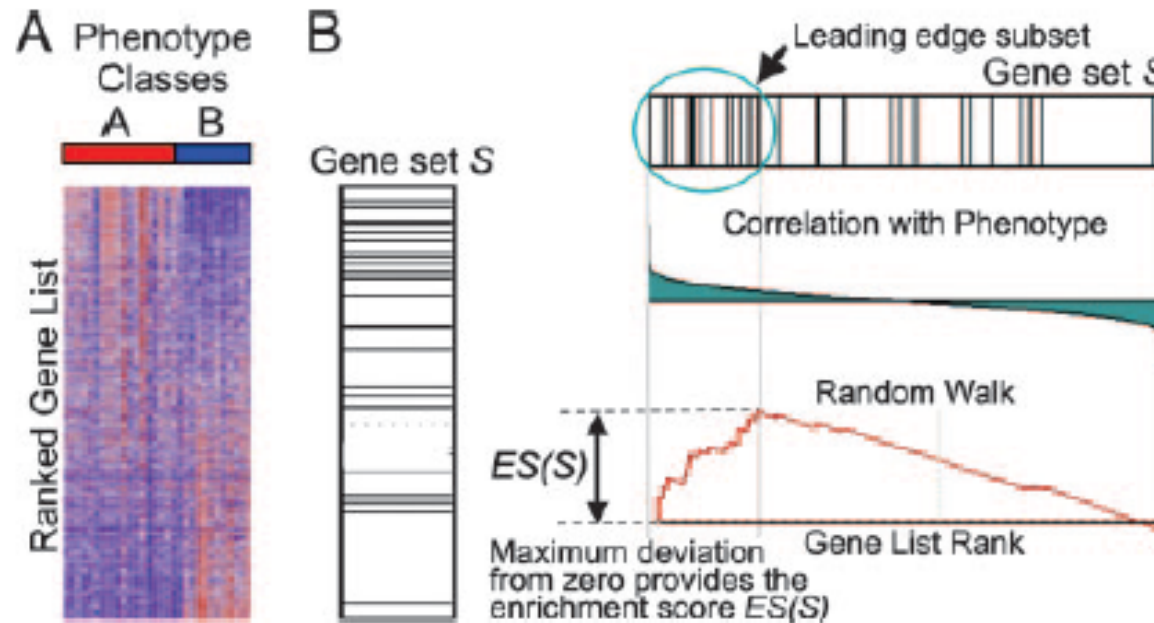
GSEA reveals biological function of ncRNAs



Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals

Mitchell Guttman^{1,2}, Ido Amit¹, Manuel Garber¹, Courtney French¹, Michael F. Lin¹, David Feldser³, Maitte Huarte^{1,6}, Or Zuk¹, Bryce W. Carey^{7,8}, John P. Cassidy^{2,8}, Moran N. Cabili⁷, Rudolf Jaenisch^{2,8}, Tarjei S. Mikkelsen^{1,4}, Tyler Jacks^{1,9}, Nir Hacohen^{1,9}, Bradley E. Bernstein^{1,10,11}, Manolis Kellis^{1,2}, Aviv Regev^{1,2}, John L. Rinn^{1,6,11} & Eric S. Lander^{1,2,7,*}

Excursion: Gene set enrichment analysis



Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles

Aravind Subramanian^{2b}, Pablo Tamayo^{2b}, Vamsi K. Mootha^{2c}, Sayan Mukherjee^d, Benjamin L. Ebert^{2e}, Michael A. Gillette^{2f}, Amanda Paulovich^{2g}, Scott L. Pomeroy^h, Todd R. Golub^{2a*}, Eric S. Lander^{2c,1,1k}, and Jill P. Mesirov^{2,k}

Summary of functional annotation

- strategy for uncovering mechanisms: generalization of examples
 - diverse molecular mechanism: antisense, association with complexes
- strategy for uncovering biological function: association with known proteins
 - ncRNAs are expressed very tissue-specifically

Conclusions

- **Many different technologies exist for identifying ncRNA**
 - based on the method you find different ncRNA types
 - be aware when comparing different data sets!
- **Computational data mining of published data has been successfully applied to identify novel ncRNAs**
 - generally identify more transcripts than single studies
 - less sequence conservation than proteins makes it more challenging to identify functional ncRNAs (use positional conservation)
 - be aware that many ncRNAs are tissue-specific!
- **Functional insights are starting to emerge**
 - generalizing known examples has been successful
 - associations of ncRNAs and genes can give insights into biological functions
 - nevertheless, we are still far away from understanding them as we do proteins